

# Bridging the Divide: Reconsidering Softmax and Linear Attention

Dongchen Han\* Yifan Pu\* Zhuofan Xia\* Yizeng Han Xuran Pan

Xiu Li Jiwen Lu Shiji Song Gao Huang



清华大学  
Tsinghua University



# Background



Transformers has a **Quadratic Complexity**  $\mathcal{O}(N^2 d)$  *with respect to sequence length.*

## High Resolution Images

## Videos





## Softmax Attention

- ✓ High expressive capability
- × **Quadratic complexity**  $\mathcal{O}(N^2 d)$

$$S_i = \left[ \frac{\exp(Q_i^\top K_1)}{\sum_{j=1}^N \exp(Q_i^\top K_j)}, \dots, \frac{\exp(Q_i^\top K_N)}{\sum_{j=1}^N \exp(Q_i^\top K_j)} \right]^\top$$

$$O_i^S = S_i^\top V$$



## Linear Attention

- × **Inferior performance**
- ✓ Linear complexity  $\mathcal{O}(Nd^2)$

$$L_i = \left[ \frac{\phi(Q_i)^\top \phi(K_1)}{\sum_{j=1}^N \phi(Q_i)^\top \phi(K_j)}, \dots, \frac{\phi(Q_i)^\top \phi(K_N)}{\sum_{j=1}^N \phi(Q_i)^\top \phi(K_j)} \right]^\top$$

$$O_i^L = L_i^\top V = \frac{\phi(Q_i)^\top \left( \sum_{j=1}^N \phi(K_j) V_j^\top \right)}{\phi(Q_i)^\top \left( \sum_{j=1}^N \phi(K_j) \right)}$$



# Injectivity of Attention Function

**Softmax Attn:**  $S_K: \mathbb{R}^d \rightarrow \mathbb{R}^N$ ,  $S_K(Q_i) = \left[ \frac{\exp(Q_i^\top K_1)}{\sum_{j=1}^N \exp(Q_i^\top K_j)}, \dots, \frac{\exp(Q_i^\top K_N)}{\sum_{j=1}^N \exp(Q_i^\top K_j)} \right]^\top$

**Linear Attn:**  $L_K: \mathbb{R}^d \rightarrow \mathbb{R}^N$ ,  $L_K(Q_i) = \left[ \frac{\phi(Q_i)^\top \phi(K_1)}{\sum_{j=1}^N \phi(Q_i)^\top \phi(K_j)}, \dots, \frac{\phi(Q_i)^\top \phi(K_N)}{\sum_{j=1}^N \phi(Q_i)^\top \phi(K_j)} \right]^\top$

## Proposition 1 -- Softmax Attn is **Injective**:

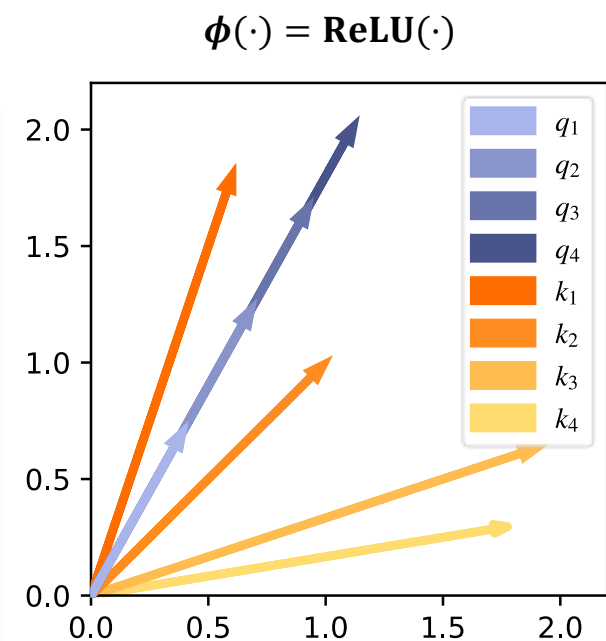
Given  $K \in \mathbb{R}^{N \times d}$  with  $\text{rank}(K) = d$ ,  $\text{rank}([K, 1]) = d + 1$ .  $\forall p, q \in \mathbb{R}^d, p \neq q$ , we have  $S_K(p) \neq S_K(q)$ .

## Proposition 2 -- Linear Attn is **Not Injective**:

Let  $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a continuous function.  $\exists p, q \in \mathbb{R}^d, p \neq q$ , s. t.  $L_K(p) = L_K(q)$ .

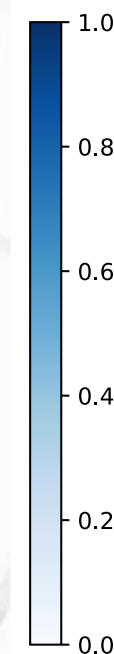
# Injectivity of Attention Function

Code



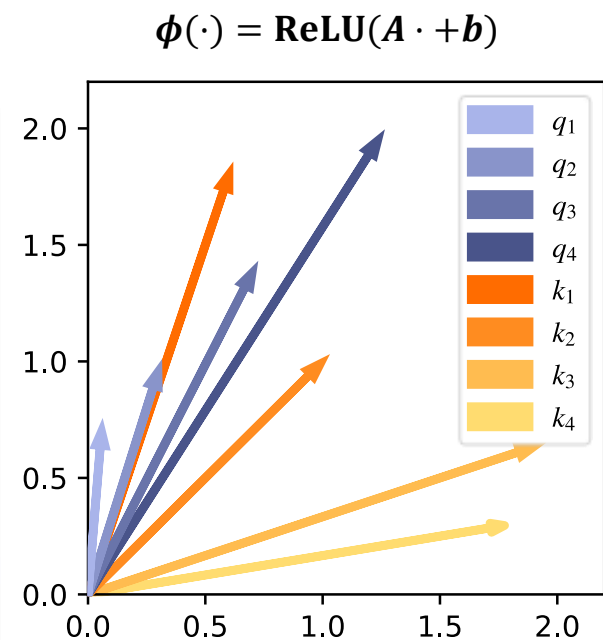
	$\text{Attn}(q_1)$	$\text{Attn}(q_2)$	$\text{Attn}(q_3)$	$\text{Attn}(q_4)$																
Softmax Attention	<table border="1"><tr><td>0.31</td><td>0.23</td></tr><tr><td>0.25</td><td>0.20</td></tr></table>	0.31	0.23	0.25	0.20	<table border="1"><tr><td>0.36</td><td>0.22</td></tr><tr><td>0.25</td><td>0.17</td></tr></table>	0.36	0.22	0.25	0.17	<table border="1"><tr><td>0.41</td><td>0.20</td></tr><tr><td>0.24</td><td>0.15</td></tr></table>	0.41	0.20	0.24	0.15	<table border="1"><tr><td>0.44</td><td>0.19</td></tr><tr><td>0.23</td><td>0.13</td></tr></table>	0.44	0.19	0.23	0.13
0.31	0.23																			
0.25	0.20																			
0.36	0.22																			
0.25	0.17																			
0.41	0.20																			
0.24	0.15																			
0.44	0.19																			
0.23	0.13																			
Linear Attention	<table border="1"><tr><td>0.32</td><td>0.23</td></tr><tr><td>0.25</td><td>0.19</td></tr></table>	0.32	0.23	0.25	0.19	<table border="1"><tr><td>0.32</td><td>0.23</td></tr><tr><td>0.25</td><td>0.19</td></tr></table>	0.32	0.23	0.25	0.19	<table border="1"><tr><td>0.32</td><td>0.23</td></tr><tr><td>0.25</td><td>0.19</td></tr></table>	0.32	0.23	0.25	0.19	<table border="1"><tr><td>0.32</td><td>0.23</td></tr><tr><td>0.25</td><td>0.19</td></tr></table>	0.32	0.23	0.25	0.19
0.32	0.23																			
0.25	0.19																			
0.32	0.23																			
0.25	0.19																			
0.32	0.23																			
0.25	0.19																			
0.32	0.23																			
0.25	0.19																			
InLine Attention (Ours)	<table border="1"><tr><td>0.31</td><td>0.23</td></tr><tr><td>0.25</td><td>0.20</td></tr></table>	0.31	0.23	0.25	0.20	<table border="1"><tr><td>0.35</td><td>0.22</td></tr><tr><td>0.26</td><td>0.17</td></tr></table>	0.35	0.22	0.26	0.17	<table border="1"><tr><td>0.39</td><td>0.21</td></tr><tr><td>0.26</td><td>0.14</td></tr></table>	0.39	0.21	0.26	0.14	<table border="1"><tr><td>0.42</td><td>0.21</td></tr><tr><td>0.26</td><td>0.12</td></tr></table>	0.42	0.21	0.26	0.12
0.31	0.23																			
0.25	0.20																			
0.35	0.22																			
0.26	0.17																			
0.39	0.21																			
0.26	0.14																			
0.42	0.21																			
0.26	0.12																			

$k_1$	$k_2$
$k_3$	$k_4$



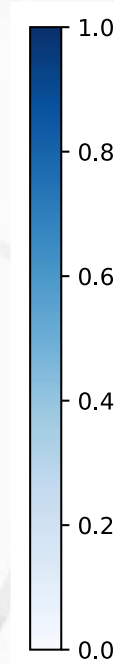
# Injectivity of Attention Function

Code



	$\text{Attn}(q_1)$	$\text{Attn}(q_2)$	$\text{Attn}(q_3)$	$\text{Attn}(q_4)$	
<b>Softmax Attention</b>	$k_1$	0.36	0.37	0.39	0.41
	$k_2$	0.25	0.23	0.21	0.19
<b>Linear Attention</b>	$k_3$	0.21	0.22	0.24	0.25
	$k_4$	0.18	0.17	0.16	0.14
<b>InLine Attention (Ours)</b>	$k_1$	0.31	0.31	0.31	0.31
	$k_2$	0.24	0.24	0.24	0.24
<b>InLine Attention (Ours)</b>	$k_3$	0.25	0.25	0.25	0.25
	$k_4$	0.21	0.21	0.21	0.21
<b>InLine Attention (Ours)</b>	$k_1$	0.35	0.37	0.40	0.45
	$k_2$	0.23	0.23	0.23	0.22
<b>InLine Attention (Ours)</b>	$k_3$	0.24	0.24	0.24	0.24
	$k_4$	0.18	0.16	0.13	0.10

$k_1$	$k_2$
$k_3$	$k_4$



# Injective Linear Attention



**Injective Linear Attn:**

$$\text{InL}_K(Q_i) = [\phi(Q_i)^\top \phi(K_1), \dots, \phi(Q_i)^\top \phi(K_N)]^\top - \frac{1}{N} \sum_{j=1}^N \phi(Q_i)^\top \phi(K_j) + \frac{1}{N}$$

**Proposition 3 -- Injective Linear Attn is *Injective*:**

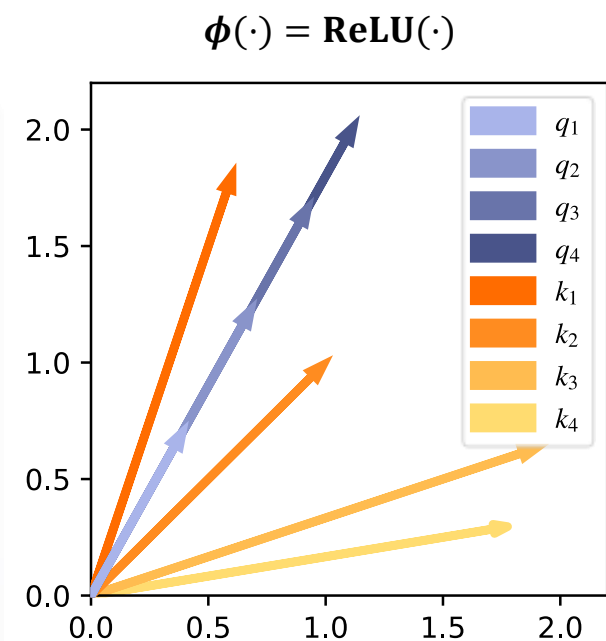
*Let  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^d$  be an injective map.*

*Given  $K \in \mathbb{R}^{N \times d}$  with  $\text{rank}(\phi(K)) = d, \text{rank}([\phi(K), 1]) = d + 1$ .*

*$\forall p, q \in \mathbb{R}^d, p \neq q$ , we have  $\text{InL}_K(p) \neq \text{InL}_K(q)$ .*

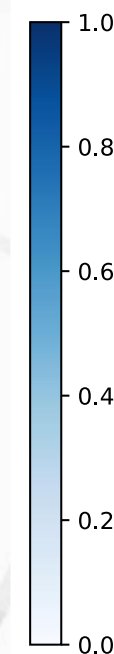
# Injectivity of Attention Function

Code



	$\text{Attn}(q_1)$	$\text{Attn}(q_2)$	$\text{Attn}(q_3)$	$\text{Attn}(q_4)$																
Softmax Attention	<table border="1"><tr><td>0.31</td><td>0.23</td></tr><tr><td>0.25</td><td>0.20</td></tr></table>	0.31	0.23	0.25	0.20	<table border="1"><tr><td>0.36</td><td>0.22</td></tr><tr><td>0.25</td><td>0.17</td></tr></table>	0.36	0.22	0.25	0.17	<table border="1"><tr><td>0.41</td><td>0.20</td></tr><tr><td>0.24</td><td>0.15</td></tr></table>	0.41	0.20	0.24	0.15	<table border="1"><tr><td>0.44</td><td>0.19</td></tr><tr><td>0.23</td><td>0.13</td></tr></table>	0.44	0.19	0.23	0.13
0.31	0.23																			
0.25	0.20																			
0.36	0.22																			
0.25	0.17																			
0.41	0.20																			
0.24	0.15																			
0.44	0.19																			
0.23	0.13																			
Linear Attention	<table border="1"><tr><td>0.32</td><td>0.23</td></tr><tr><td>0.25</td><td>0.19</td></tr></table>	0.32	0.23	0.25	0.19	<table border="1"><tr><td>0.32</td><td>0.23</td></tr><tr><td>0.25</td><td>0.19</td></tr></table>	0.32	0.23	0.25	0.19	<table border="1"><tr><td>0.32</td><td>0.23</td></tr><tr><td>0.25</td><td>0.19</td></tr></table>	0.32	0.23	0.25	0.19	<table border="1"><tr><td>0.32</td><td>0.23</td></tr><tr><td>0.25</td><td>0.19</td></tr></table>	0.32	0.23	0.25	0.19
0.32	0.23																			
0.25	0.19																			
0.32	0.23																			
0.25	0.19																			
0.32	0.23																			
0.25	0.19																			
0.32	0.23																			
0.25	0.19																			
InLine Attention (Ours)	<table border="1"><tr><td>0.31</td><td>0.23</td></tr><tr><td>0.25</td><td>0.20</td></tr></table>	0.31	0.23	0.25	0.20	<table border="1"><tr><td>0.35</td><td>0.22</td></tr><tr><td>0.26</td><td>0.17</td></tr></table>	0.35	0.22	0.26	0.17	<table border="1"><tr><td>0.39</td><td>0.21</td></tr><tr><td>0.26</td><td>0.14</td></tr></table>	0.39	0.21	0.26	0.14	<table border="1"><tr><td>0.42</td><td>0.21</td></tr><tr><td>0.26</td><td>0.12</td></tr></table>	0.42	0.21	0.26	0.12
0.31	0.23																			
0.25	0.20																			
0.35	0.22																			
0.26	0.17																			
0.39	0.21																			
0.26	0.14																			
0.42	0.21																			
0.26	0.12																			

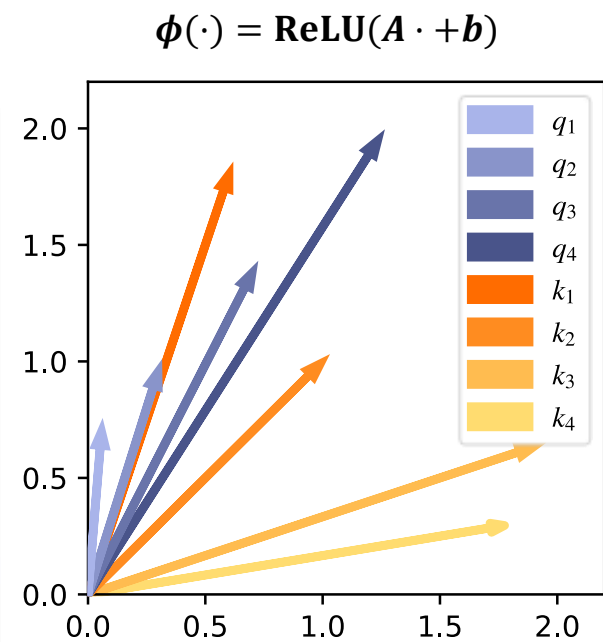
$k_1$	$k_2$
$k_3$	$k_4$





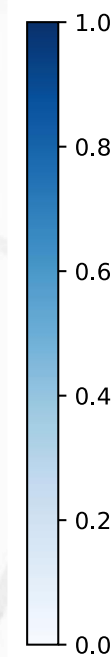
# Injectivity of Attention Function

Code



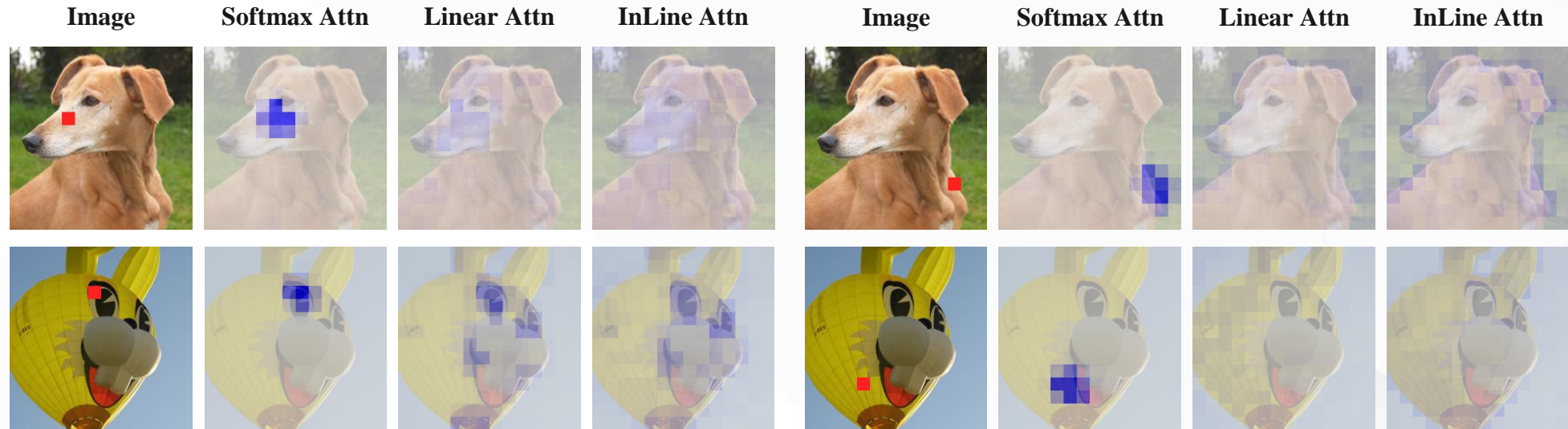
	$\text{Attn}(q_1)$	$\text{Attn}(q_2)$	$\text{Attn}(q_3)$	$\text{Attn}(q_4)$
<b>Softmax Attention</b>	0.36	0.25	0.37	0.23
	0.21	0.18	0.22	0.17
<b>Linear Attention</b>	0.31	0.24	0.31	0.24
	0.25	0.21	0.25	0.21
<b>InLine Attention (Ours)</b>	0.35	0.23	0.37	0.23
	0.24	0.18	0.24	0.16

$k_1$	$k_2$
$k_3$	$k_4$



# Local Modeling Capability

Code



Mask Out Position	None	Loc. 3×3	Loc. 5×5	Loc. 7×7	Rand 9	Rand 25	Rand 49
Softmax Attn	72.2	51.6	24.3	9.0	71.7	71.5	71.1
InLine Attn	70.0	58.0	40.0	20.0	70.0	69.9	69.5

# Local Modeling Capability



## Injective Linear Attn:

$$\text{InL}_K(Q_i) = [\phi(Q_i)^\top \phi(K_1), \dots, \phi(Q_i)^\top \phi(K_N)]^\top - \frac{1}{N} \sum_{j=1}^N \phi(Q_i)^\top \phi(K_j) + \frac{1}{N}$$

## InLine Attention Module:

$$O_i = \text{InL}_K(Q_i)^\top V + \sum_{j=1}^9 r_j V_j^{N(i)}, \quad r = \text{MLP}(\bar{x})$$



✓ The impact of injective property:

Kernel Function $\phi(\cdot)$	ReLU( $\cdot$ )	ReLU( $A \cdot +b$ )	LeakyReLU( $\cdot$ )	Identity( $\cdot$ )
Linear Attn	77.3	70.2	1.5	0.2
InLine Attn	79.8	80.0	79.8	80.2

✓ Local modeling ability:

	Window	FLOPs	#Param	Acc.
InLine-Swin-T w/o res.	$7^2$	4.5G	30M	80.3
	$14^2$	4.5G	30M	80.4
	$28^2$	4.5G	30M	80.2
	$56^2$	4.5G	30M	80.2

	Window	FLOPs	#Param	Acc.
InLine-Swin-T w/ res.	$7^2$	4.5G	30M	81.6
	$14^2$	4.5G	30M	82.1
	$28^2$	4.5G	30M	82.3
	$56^2$	4.5G	30M	82.4



✓ Performances on ImageNet-1K:

Method	Reso	#Params	FLOPs	Top-1
DeiT-T [30]	224 <sup>2</sup>	5.7M	1.2G	72.2
<b>InLine-DeiT-T</b>	224 <sup>2</sup>	6.5M	1.1G	<b>74.5 (+2.3)</b>
DeiT-B	224 <sup>2</sup>	86.6M	17.6G	81.8
<b>InLine-DeiT-B</b>	448 <sup>2</sup>	23.8M	17.2G	<b>82.3 (+0.5)</b>
PVT-S	224 <sup>2</sup>	24.5M	3.8G	79.8
<b>InLine-PVT-S</b>	224 <sup>2</sup>	21.6M	3.9G	<b>82.0 (+2.2)</b>
PVT-L	224 <sup>2</sup>	61.4M	9.8G	81.7
<b>InLine-PVT-L</b>	224 <sup>2</sup>	50.2M	10.2G	<b>83.6 (+1.9)</b>

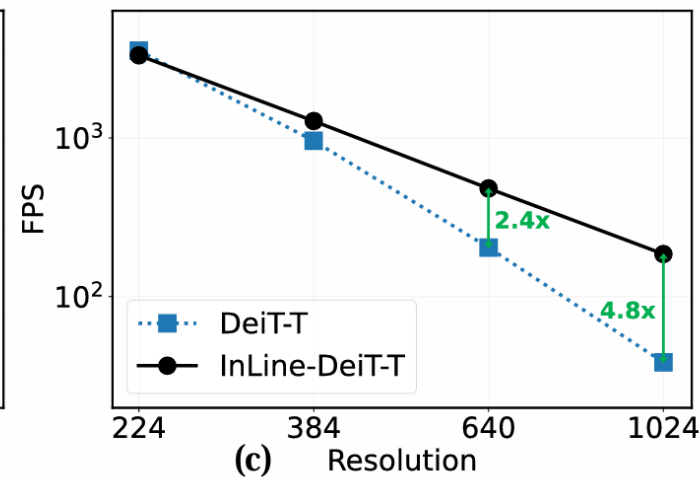
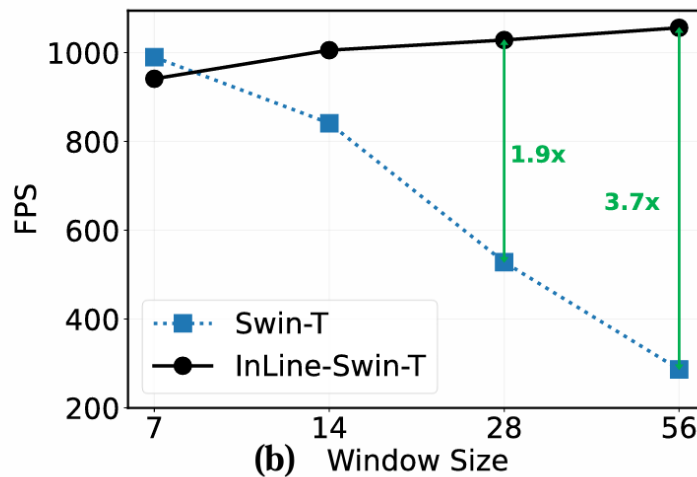
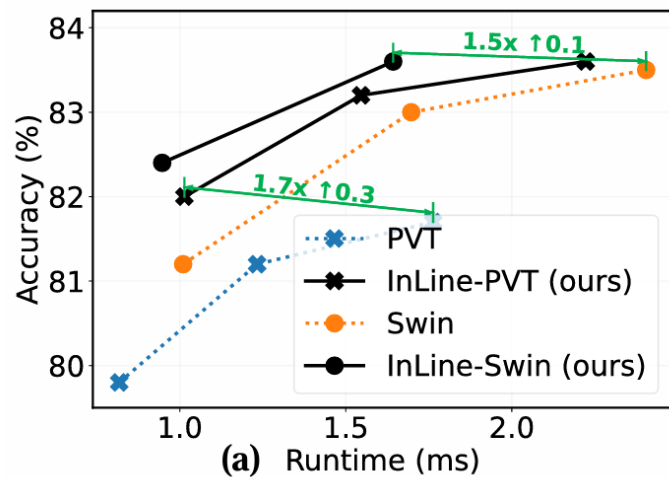
Method	Reso	#Params	FLOPs	Top-1
Swin-T [19]	224 <sup>2</sup>	29M	4.5G	81.3
<b>InLine-Swin-T</b>	224 <sup>2</sup>	30M	4.5G	<b>82.4 (+1.1)</b>
Swin-S	224 <sup>2</sup>	50M	8.7G	83.0
<b>InLine-Swin-S</b>	224 <sup>2</sup>	50M	8.7G	<b>83.6 (+0.6)</b>
Swin-B	224 <sup>2</sup>	88M	15.4G	83.5
<b>InLine-Swin-B</b>	224 <sup>2</sup>	88M	15.4G	<b>84.1 (+0.6)</b>
Swin-B	384 <sup>2</sup>	88M	47.0G	84.5
<b>InLine-Swin-B</b>	384 <sup>2</sup>	88M	45.2G	<b>85.0 (+0.5)</b>

# Empirical Study

Code



✓ Speed measurements:





✓ Performances on downstream tasks:

<b>(a) Mask R-CNN Object Detection on COCO</b>								
Method	FLOPs	Sch.	$AP^b$	$AP_{50}^b$	$AP_{75}^b$	$AP^m$	$AP_{50}^m$	$AP_{75}^m$
PVT-T	240G	1x	36.7	59.2	39.3	35.1	56.7	37.3
InLine-PVT-T	211G	1x	40.2	62.7	43.8	37.7	59.7	40.4
PVT-S	305G	1x	40.4	62.9	43.8	37.8	60.1	40.3
InLine-PVT-S	250G	1x	43.4	66.4	47.1	40.1	63.1	43.3
PVT-M	392G	1x	42.0	64.4	45.6	39.0	61.6	42.1
InLine-PVT-M	310G	1x	44.0	66.4	48.0	40.3	63.4	43.5
PVT-L	494G	1x	42.9	65.0	46.6	39.5	61.9	42.5
InLine-PVT-L	377G	1x	45.4	67.6	49.7	41.4	64.7	44.6

<b>(b) Cascade Mask R-CNN Object Detection on COCO</b>								
Method	FLOPs	Sch.	$AP^b$	$AP_{50}^b$	$AP_{75}^b$	$AP^m$	$AP_{50}^m$	$AP_{75}^m$
Swin-S	837G	3x	51.9	70.7	56.3	45.0	68.2	48.8
InLine-Swin-S	835G	3x	52.4	71.0	56.9	45.4	68.8	49.6
Swin-B	981G	3x	51.9	70.5	56.4	45.0	68.1	48.9
InLine-Swin-B	978G	3x	52.6	71.0	57.0	45.4	68.5	49.3



## Injective Linear Attention (**InLine**)

- ✓ The injectivity of attention function is of crucial importance
- ✓ Local modeling is essential to attention
- ✓ InLine: a simple, fast and effective linear attention module





清华大学  
Tsinghua University



# Thank you!

**Contact:** [hdc23@mails.tsinghua.edu.cn](mailto:hdc23@mails.tsinghua.edu.cn)

Code

